

ОБ ОДНОЙ ЗАДАЧЕ ВЕРОЯТНОСТНО-СТАТИСТИЧЕСКОГО АНАЛИЗА ВКРАПЛЕНИЙ В ДВОИЧНУЮ ЦЕПЬ МАРКОВА

В статье рассматривается (q, r) -блочная модель вкраплений в двоичную цепь Маркова, возникающая в задачах стеганографической защиты информации. Построены статистические оценки параметров модели на основе частотных статистик. Представлены результаты компьютерных экспериментов.

Ключевые слова: цепь Маркова; вкрапление; статистические оценки; блочная модель.

In this paper (q, r) -block mathematical model of embeddings in binary Markov chain that appear in steganography is considered. Estimators of model parameters based on frequencies statistics are constructed. The results of computer experiments are presented.

Key words: Markov chain; embedding; statistical estimators; block model.

1. Введение и математическая модель вкраплений. Актуальной задачей стеганографической защиты информации является задача обнаружения вкраплений [1–8], т. е. обнаружение встраивания сообщения в контейнер. В качестве контейнеров могут использоваться изображения, аудио- и видео-последовательности. В литературе из этой области недостаточно внимания уделено вероятностно-статистическим вопросам.

Цель данной статьи – разработка методов статистического анализа математической модели вкраплений. Предлагается новая (q, r) -блочная модель вкраплений в однородную цепь Маркова, построены статистические оценки параметров ε, δ для (q, r) -блочной модели при $q = 2, r = 1$. Представлены результаты компьютерных экспериментов.

Введем обозначения: $V = \{0, 1\}$ – двоичный алфавит; V_t – пространство двоичных t -мерных векторов; N – множество натуральных чисел; $I\{A\}$ – индикатор события A ; $u_{t_1}^{t_2} = (u_{t_1}, \dots, u_{t_2}) \in V_{t_2-t_1+1}$ ($t_1, t_2 \in N, t_1 \leq t_2$) – двоичная строка из $t_2 - t_1 + 1$ символов; $w(\cdot)$ – вес Хемминга; $L\{\xi\}$ – закон распределения вероятностей случайной величины ξ ; $B(\theta)$ – закон распределения вероятностей Бернулли с параметром $\theta \in [0, 1]$.

Будем предполагать, что контейнер для встраивания сообщения есть двоичная последовательность $x_1^T \in V_T, x_t \in V, t = 1, \dots, T$, длиной T , являющаяся однородной двоичной цепью Маркова 1-го порядка с симметричной матрицей вероятностей одношаговых переходов $P = (p_{j_0, j_1}), j_0, j_1 \in V$:

$$P = P(\varepsilon) = \frac{1}{2} \begin{pmatrix} 1 + \varepsilon & 1 - \varepsilon \\ 1 - \varepsilon & 1 + \varepsilon \end{pmatrix}, p_{j_0, j_1} = P\{x_{t+1} = j_1 | x_t = j_0\} = \frac{1}{2}(1 + (-1)^{j_0 + j_1} \varepsilon), |\varepsilon| < 1. \quad (1)$$

Здесь ε – параметр модели: случай $\varepsilon = 0$ соответствует схеме независимых испытаний и исследований в [7]; случай $\varepsilon > 0$ учитывает зависимость типа притяжения; $\varepsilon < 0$ – зависимость типа отталкивания. Отметим, что цепь Маркова (1) удовлетворяет условиям эргодичности [9] и имеет равномерное стационарное распределение вероятностей $\pi = (1/2, 1/2)$. Далее будем полагать, что цепь Маркова (1) является стационарной.

Обычно (см. [5]) на практике сообщение перед встраиванием в контейнер подвергается криптографическому преобразованию, устраняющему статистическую избыточность, поэтому далее полагаем, что сообщение $\xi_1^M \in V_M, M \leq T$, является последовательностью M независимых случайных величин Бернулли: $L\{\xi_t\} = B(\theta), P\{\xi_t = j\} = \theta_j, j \in V, \theta_1 = 1 - \theta_0, t = 1, \dots, M$. Как правило (см. [5]), $\{\xi_t\}$ имеет симметричное распределение вероятностей: $\theta_1 = \theta_0 = 1/2$.

Стегослуж $\gamma_1^T \in V_T$ определяет моменты времени, в которые биты сообщения ξ_1^M вкрапляются в последовательность x_1^T . Введем специальную (q, r) -блочную модель стегослужа ($q, r \in N, r \leq q$), предполагая, что длина последовательности x_1^T кратна q : $T = Kq$. Для этого вначале разобьем последовательность x_1^T на блоки длиной q : $x_{(1)} = x_1^q, x_{(2)} = x_{q+1}^q, \dots, x_{(K)} = x_{T-q+1}^q$. Введем вспомогательные независимые случайные величины $\varsigma_k \in V, L\{\varsigma_k\} = B(\delta), k = 1, \dots, K$, которые отвечают за выбор блоков $\{x_{(k)}\}$ для вкрапления сообщения: если $\varsigma_k = 1$, то в r случайно выбранных бит блока $\{x_{(k)}\}$ вкрапляются r бит сообщения; если $\varsigma_k = 0$, то вкрапление в блок $\{x_{(k)}\}$ не производится.

В результате стегослуж γ_1^T состоит из независимых блоков, имеющих следующие распределения вероятностей:

$$P\{\gamma_{(k-1)q+1}^{kq} = u_1^q\} = \begin{cases} 1 - \delta, & w(u_1^q) = 0, \\ \delta / C_q^r, & w(u_1^q) = r, \\ 0, & w(u_1^q) \notin \{0, r\}, \end{cases} \quad k = 1, \dots, K, u_1^q \in V_q. \quad (2)$$

Отметим, что для такой модели стежоключа максимальная пропускная способность стегосистемы уменьшается до $Tr/q = Kr$ бит, а мощность множества всевозможных стежоключей сжимается до $(1 + C_q^r)^K$.

Случайная стежокпоследовательность $Y_1^T \in V_T$ порождается следующим образом:

$$Y_t = \gamma_t \xi_{\gamma_t} + (1 - \gamma_t)x_t = \begin{cases} x_t, & \gamma_t = 0, \\ \xi_{\gamma_t}, & \gamma_t = 1. \end{cases} \quad (3)$$

Случайные последовательности $\{x_t\}$, $\{\xi_t\}$, $\{\gamma_t\}$ предполагаются независимыми в совокупности.

Заметим, что с практической точки зрения наибольшего внимания в рамках рассматриваемой здесь марковской модели вкраплений (3) заслуживает наиболее трудный для обнаружения вкраплений случай, когда $\theta_1 = \theta_0 = 1/2$, так как в этом случае при вкраплении одномерное распределение вероятностей не искажается: $R\{Y_t = j\} = R\{x_t = j\} = 1/2$, $j \in V$.

2. Статистическое оценивание параметров модели вкраплений

Аналогично [8] разобьем множество V_t двоичных t -мерных векторов на $t+1$ непересекающихся подмножеств:

$$V_t = \Gamma_0^{(t)} \cup \Gamma_1^{(t)} \cup \dots \cup \Gamma_t^{(t)}, \quad (4)$$

где

$$\begin{aligned} \Gamma_0^{(t)} &= \{u_1^t \in V_t : u_t = 1\}, \\ \Gamma_1^{(t)} &= \{u_1^t \in V_t : u_t = u_{t-1} = 0\}, \\ \Gamma_j^{(t)} &= \{u_1^t \in V_t : u_t = 0, u_{t-1} = \dots = u_{t-j} = 1, u_{t-j-1} = 0\}, \quad 1 < j < t, \\ \Gamma_t^{(t)} &= \{u_1^t \in V_t : u_t = 0, u_{t-1} = \dots = u_1 = 1\}. \end{aligned} \quad (5)$$

Такое разбиение порождает разбиение всевозможных траекторий фрагментов стежоключа $\gamma_1^t = u_1^t \in V_t$. Определим функции двоичных переменных $u_1^t \in V_t$, $y_1^t \in V_t$:

$$\varphi_t(u_1^t, y_1^t) = \begin{cases} \theta_{y_t}, & u_1^t \in \Gamma_0^{(t)}, \\ (1 + (-1)^{y_t + y_{t-j}} \varepsilon^j) / 2, & u_1^t \in \Gamma_j^{(t)}, 1 \leq j \leq t-1, \\ 1/2, & u_1^t \in \Gamma_t^{(t)}. \end{cases} \quad (6)$$

С учетом (6) и введенных обозначений функция правдоподобия для модели вкраплений (3) имеет вид

$$L(\varepsilon, \delta) = \sum_{u_1^T \in V_T} I \left\{ \sum_{h \in \{0, r\}} b_h(u_1^T) = 0 \right\} (1 - \delta)^{b_0(u_1^T)} (\delta / C_q^r)^{b_r(u_1^T)} \prod_{t=1}^T \varphi_t(u_1^t, y_1^t), \quad (7)$$

где $y_1^T \in V_T$ – наблюдаемая стежокпоследовательность; $b_h(u_1^T) = \sum_{k=1}^{T/q} I \{w(u_{q(k-1)+1}^{qk}) = h\}$, $h \in \{0, \dots, q\}$. Согласно определению (7) вычисление одного значения функции правдоподобия $L(\varepsilon, \delta)$ при фиксированных параметрах ε, δ имеет вычислительную сложность порядка $O(T(1 + C_q^r)^{T/q})$.

Оценки максимального правдоподобия $\hat{\varepsilon}, \hat{\delta}$ параметров ε, δ определяются как решение экстремальной задачи: $L(\varepsilon, \delta) \rightarrow \max_{\varepsilon \in (-1, 1), \delta \in [0, 1]}$. Решение этой задачи возможно только численными методами (например, методом табулирования $L(\varepsilon, \delta)$ на сетке или методом градиентного спуска), требующими вычисления функции правдоподобия для заданной последовательности точек.

Лемма 1. Если имеет место модель вкраплений (3), то при $q > r$, $t > 2r + 1$

$$P \left\{ \gamma_1^t \in \bigcup_{j=2r+2}^t \Gamma_j^{(t)} \right\} = 0, \quad t > 2r + 1, \quad (8)$$

$$\begin{aligned} P \{ Y_t = y_t \mid Y_1^{t-1} = y_1^{t-1}, \gamma_1^t = u_1^t \} &= \lambda_t(u_{t-2r-1}^t, y_{t-2r-1}^t) = \\ &= \begin{cases} \theta_{y_t}, & u_1^t \in \Gamma_0^{(t)}, \\ (1 + (-1)^{y_t + y_{t-j}} \varepsilon^j) / 2, & u_1^t \in \Gamma_j^{(t)}, 1 \leq j \leq 2r + 1, \end{cases} \end{aligned} \quad (9)$$

и стежокпоследовательность $\{Y_t\}$, определяемая в (3), при фиксированной последовательности $\{\gamma_t\}$ является управляемой цепью Маркова условного порядка $s_t \in \{0, \dots, 2r + 1\}$, причем порядок s_t зависит от ключевой последовательности $\{\gamma_t\}$: $s_t = j$, если $u_1^t \in \Gamma_j^{(t)}$.

Доказательство. Соотношение (8) вытекает из (2) и обозначений (4), (5). Для доказательства (9) достаточно проверить марковское свойство последовательности Y_t при фиксированной управляющей последовательности γ_t .

Лемма 1 позволяет построить полиномиальный относительно T алгоритм вычисления значения функции правдоподобия $L(\varepsilon, \delta)$ для (q, r) -блочной модели вкраплений при $q > r$.

При $q = 2, r = 1$ и $\theta_0 = \theta_1 = 1/2$ удобно использовать теорему 1 для построения частотных оценок параметров ε, δ .

Теорема 1. Если $q = 2, r = 1$ и $\theta_0 = \theta_1 = 1/2$, то матрица усредненных по γ_{t-3}^t ($t > 2r + 1$) вероятностей одношаговых переходов управляемой цепи Маркова 3-го порядка имеет вид

$$\bar{P}_t = (\bar{p}_{y_{t-3}, y_{t-2}, y_{t-1}, y_t}) = \begin{cases} (CA_{2k} : 1_8 - CA_{2k}), t = 2k, k \in N, \\ (CA_{2k-1} : 1_8 - CA_{2k-1}), t = 2k-1, k \in N, \end{cases} \quad (10)$$

где $1_8 = (1, 1, 1, 1, 1, 1, 1, 1)'$,

$$A_{2k} = \begin{pmatrix} 1/2 \\ \varepsilon(1-\delta)/2 \\ \varepsilon^2\delta(1-\delta/2)/4 \\ \varepsilon^3\delta^2/8 \end{pmatrix}, \quad A_{2k-1} = \begin{pmatrix} 1/2 \\ \varepsilon(1-\delta+\delta^2/4)/2 \\ \varepsilon^2\delta(1-\delta/2)/4 \\ 0 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & 1 \\ 1 & 1 & -1 & 1 \\ 1 & -1 & -1 & 1 \\ 1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & -1 \end{pmatrix}.$$

Доказательство. С учетом леммы 1 имеем

$$\bar{p}_{y_{t-3}, y_{t-2}, y_{t-1}, y_t} = P\{Y_t = y_t \mid Y_{t-3}^{t-1} = y_{t-3}^{t-1}\} = \sum_{j=0}^3 P\{\gamma_1^t \in \Gamma_j^{(t)}\} \lambda_t(u_{t-3}^t, y_{t-3}^t).$$

Подставляя сюда следующие из (2) значения вероятностей:

$$P\{\gamma_1^{2k} \in \Gamma_j^{(2k)}\} = \begin{cases} \delta/2, j=0, \\ 1-\delta, j=1, \\ \delta(1-\delta/2)/2, j=2, \\ \delta^2/4, j=3, \end{cases} \quad P\{\gamma_1^{2k-1} \in \Gamma_j^{(2k-1)}\} = \begin{cases} \delta/2, j=0, \\ 1-\delta+\delta^2/4, j=1, \\ \delta(1-\delta/2)/2, j=2, \\ 0, j=3, \end{cases}$$

получаем (10).

Обозначим: $J_{t_0}^{(0)} = \{t \in [t_0, T] : t = 2k, k \in N\}$, $J_{t_0}^{(1)} = \{t \in [t_0, T] : t = 2k-1, k \in N\}$ – множества четных и нечетных моментов времени соответственно, начиная от заданного начального момента времени $t_0 \in N$; $f_{v_0 v_1 v_2 v_3}^{(j)} = \sum_{t \in J_{t_0}^{(j)}} I\{Y_{t-3}^t = v_0^3\}$, $j \in V$, $v_0^3 \in V_4$, – абсолютная частота 4-граммы для четных моментов времени ($j=0$) и нечетных ($j=1$) соответственно. Частотные оценки вероятностей (10) по наблюдаемой последовательности $y_1^T \in V_T$ имеют вид

$$\hat{p}_{v_0 v_1 v_2 v_3}^{(j)} = \frac{f_{v_0 v_1 v_2 v_3}^{(j)}}{f_{v_0 v_1 v_2 0}^{(j)} + f_{v_0 v_1 v_2 1}^{(j)}}, \quad j \in V, \quad v_0^3 \in V_4.$$

Определим линейные функции от частот $\{f_{v_0 v_1 v_2 0}^{(j)}\}$:

$$s_1^{(j)} = \hat{p}_{0000}^{(j)} + \hat{p}_{0010}^{(j)} + \hat{p}_{0100}^{(j)} + \hat{p}_{1000}^{(j)} + \hat{p}_{1100}^{(j)} - \hat{p}_{0110}^{(j)} - \hat{p}_{1010}^{(j)}, \quad s_2^{(j)} = s_1^{(j)} + 2\hat{p}_{1010}^{(j)} - 2\hat{p}_{1100}^{(j)}.$$

Из (10) для построения статистических оценок параметров модели вкраплений имеем систему уравнений:

$$\begin{cases} s_1^{(0)} = 2 + 2\varepsilon(1-\delta), \\ s_1^{(1)} = 2 + 2\varepsilon(1-\delta/2)^2, \\ s_2^{(j)} = 2 + \varepsilon^2\delta(1-\delta/2), j \in V. \end{cases} \quad (11)$$

При известном параметре ε , решая систему уравнений (11), получаем подстановочную (plug-in) оценку параметра $\delta \in [0, 1]$, являющуюся функцией частот $\{f_{v_0 v_1 v_2 v_3}^{(j)}\}$:

$$\hat{\delta}_{2k}(\varepsilon) = \begin{cases} 0, & g_0 > 1, \\ 1 - g_0, & g_0 \in [0, 1], \\ 1, & g_0 < 0, \end{cases} \quad \hat{\delta}_{2k-1}(\varepsilon) = \begin{cases} 0, & g_1 > 1, \\ 2 - 2\sqrt{g_1}, & g_1 \in [1/4, 1], \\ 1, & g_1 < 1/4. \end{cases} \quad g_j = \frac{s_1^{(j)} - 2}{2\varepsilon}, \quad j \in V, \quad (12)$$

Усредняя статистические оценки для параметра δ при четных и нечетных моментах времени, получаем еще одну статистическую оценку:

$$\hat{\delta}(\varepsilon) = (\hat{\delta}_{2k}(\varepsilon) + \hat{\delta}_{2k-1}(\varepsilon)) / 2. \quad (13)$$

Если неизвестны оба параметра ε, δ модели вкраплений, то для четных моментов времени статистическая оценка $\hat{\delta}_{2k} \in [0, 1]$ имеет вид

$$\hat{\delta}_{2k} = \begin{cases} 0, & g > 1, \\ 1 - \sqrt{g}, & g \in [0, 1], \\ 1, & g < 0. \end{cases} \quad g = g(s_1^{(0)}, s_2^{(0)}) = \frac{(s_1^{(0)} - 2)^2}{(s_1^{(0)} - 2)^2 + 8(s_2^{(0)} - 2)^2}, \quad (14)$$

Для нечетных моментов времени статистическая оценка $\hat{\delta}_{2k-1} \in [0, 1]$ имеет вид

$$\hat{\delta}_{2k-1} = \min_{\alpha \in [0, 1]} |(s_2^{(1)} - 2)(2 - \alpha)^3 - 2(s_1^{(1)} - 2)^2 \alpha|. \quad (15)$$

Оценки параметра ε модели при $\hat{\delta}_{2k} < 1, \hat{\delta}_{2k-1} < 1$ строятся следующим образом: $\hat{\varepsilon}_{2k} = \frac{s_1^{(0)} - 2}{2(1 - \hat{\delta}_{2k})}$,

$$\hat{\varepsilon}_{2k-1} = \frac{s_1^{(0)} - 2}{2(1 - \hat{\delta}_{2k-1})}.$$

3. Численные результаты. Для модели вкраплений (3) при известном параметре ε методом Монте-Карло вычислены выборочные среднеквадратические ошибки $v\{\hat{\delta}(\varepsilon)\} = \frac{1}{N} \sum_{n=1}^N (\hat{\delta}^{(n)}(\varepsilon) - \delta)^2$ оценивания параметра δ ($\hat{\delta}^{(n)}(\varepsilon)$ – оценка параметра δ для n -й реализации стегопоследовательности) на основе частотных статистик (12), (13) при следующих значениях параметров: $r = 1, q = 2, \varepsilon = 0,3, \delta = 0,2$, число прогонов $N = 10^3$.

На рис. 1 изображены графики зависимостей $v\{\hat{\delta}(\varepsilon)\}$ от длины стегопоследовательности T для трех статистических оценок. Закрашенными кружками отмечены графики для частотных оценок (13), незакрашенными – для оценок $\hat{\delta}_{2k}(\varepsilon)$ по четным моментам времени, треугольниками – для оценок $\hat{\delta}_{2k-1}(\varepsilon)$.

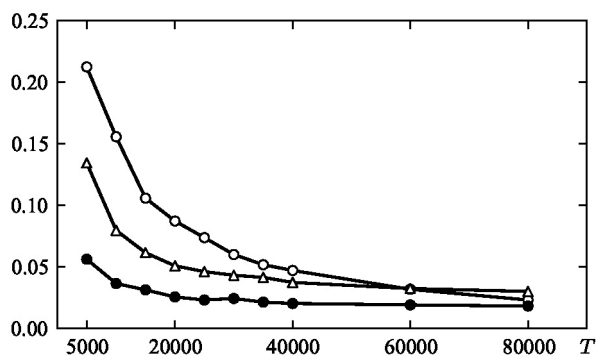


Рис. 1. Зависимость $v\{\hat{\delta}(\varepsilon)\}$ от длины T при $\varepsilon = 0,3, \delta = 0,2, N = 10^3$

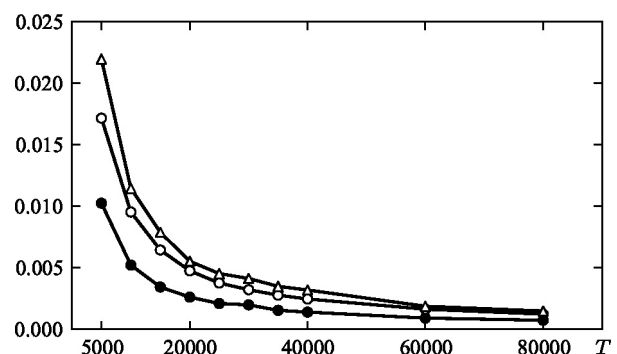


Рис. 2. Зависимость $v\{\hat{\delta}\}$ от длины T при $\varepsilon = 0,3, \delta = 0,2, N = 10^3$

При неизвестных параметрах модели методом Монте-Карло вычислены выборочные среднеквадратические ошибки оценивания параметра δ на основе (14), (15) при следующих значениях параметров $r = 1, q = 2, \varepsilon = 0,3, \delta = 0,2, N = 10^3$.

На рис. 2 изображены графики зависимостей $v\{\hat{\delta}\}$ от длины стегопоследовательности T для трех статистических оценок. Закрашенными кружками отмечены графики для усредненных оценок при четных и нечетных моментах времени, незакрашенными – для частотных оценок $\hat{\delta}_{2k}$, определенных в (14), треугольниками – для оценок $\hat{\delta}_{2k-1}$, определенных в (15).

БИБЛИОГРАФИЧЕСКИЙ СПИСОК

1. Sullivan K., Madhow U., Chandrasekaran S., Manjunath B. Steganalysis of spread spectrum data hiding exploiting cover memory // Proc. SPIE, Electronic Imaging, Security, Steganography and Watermarking of Multimedia Contents VII. San Jose, 2005. Vol. 5681. P. 38–46.
2. Pevny T., Bas P., Fridrich J. Steganalysis by subtractive pixel adjacency matrix // Proceedings of the 11th ACM Multimedia and Security Workshop. Princeton, 2009. P. 75–84.
3. Bas P., Filler T., Pevny T. Break Our Steganographic System – the ins and outs of organizing BOSS // Proceedings of Information Hiding Conference. Prague, 2011.
4. Kharin Yu. S., Vecherko E. V. On statistical hypotheses testing of embedding // Proceedings of the 9th Intern. Conf. Computer Data Analysis and Modeling. Minsk, 2010. Vol. 2. P. 26–29.
5. Харин Ю. С., Берник В. И., Матвеев Г. В., Агиевич С. В. Математические и компьютерные основы криптологии. Минск; М., 2003.
6. Харин Ю. С., Вечерко Е. В. О некоторых задачах статистической проверки гипотез в стеганографии // Весці НАН Беларусі. Сер. фіз.-мат. навук. 2010. № 4. С. 5–12.
7. Пономарев К. И. Параметрическая модель вкрапления и ее статистический анализ // Дискретная математика. 2009. Т. 21, вып. 4. С. 148–157.
8. Харин Ю. С., Вечерко Е. В. Статистическое оценивание параметров модели вкраплений в двоичную цепь Маркова // Дискретная математика. 2013. Т. 25, вып. 2. С. 135–148.
9. Кемени Дж., Снелл Дж. Конечные цепи Маркова. М., 1982.

Поступила в редакцию 12.06.13.

Егор Валентинович Вечерко – аспирант кафедры математического моделирования и анализа данных. Научный руководитель – член-корреспондент НАН Беларуси, доктор физико-математических наук, профессор, заведующий кафедрой математического моделирования и анализа данных, директор НИИ ППМИ Ю. С. Харин.